# [Ve]tting [R]esponse [I]ntegrity from cross-[T]alk dependencies in [A]dversarial [S]urveys:
## *"Who can catch a liar?"*

Nicholas Sizemore, Royce Lee, Robert Gibbons, and Ishanu Chattopadhyay[1,4,5,7]★

[1]Department of Medicine, University of Chicago, Chicago, IL 60637, USA
[4]Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL 60637, USA
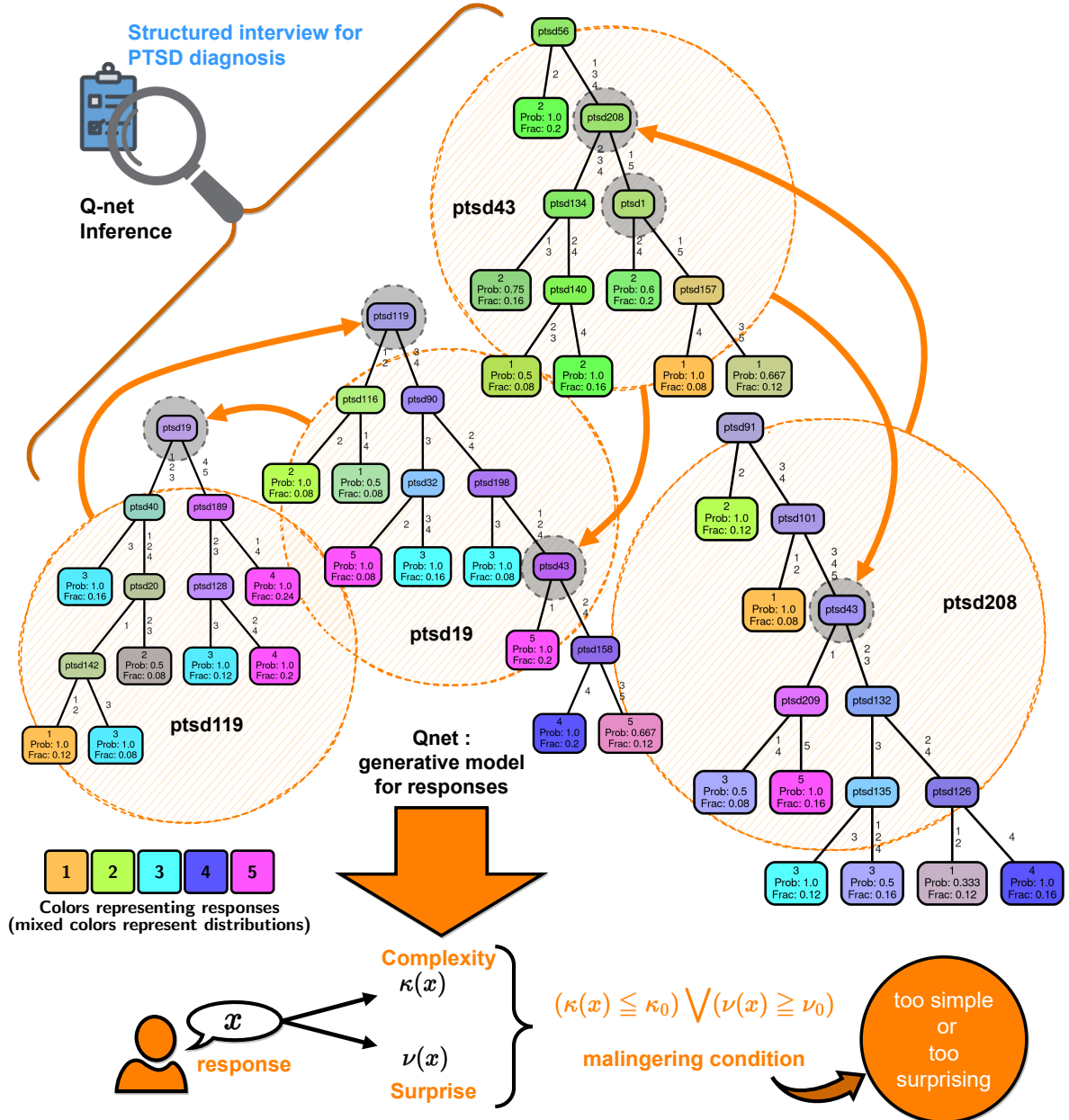[5]Committee on Genetics, Genomics & Systems Biology, University of Chicago, Chicago, IL 60637, USA

★To whom correspondence should be addressed: e-mail: ishanu@uchicago.edu.

*Abstract:* **Malingering[1,2] or faking the symptoms of a mental health disorder can confound structured diagnostic interviews and hinder clinical psychiatric assessments. We introduce an artificial intelligence (AI) framework for detecting symptom fabrication in Post-Traumatic Stress Disorder (PTSD) diagnoses, for which malingering is a known problem[3,4], partially ascribable to the potential for secondary financial gain from positive diagnoses. Algorithm VeRITAS employs novel generative AI to infer statistical dependencies inherent in true response patterns, and flags responses which violate these subtle constraints. Mimicing these emergent patterns is difficult even with psychiatric training, resulting in a robust mechanism for recognizing fabrications. With a study sample of $n = 624$ patients, VeRITAS is estimated to have a Area Under the Curve (AUC) of $\geq 0.95 \pm 0.02$, with one sensitivity/specificity pair of 86% and 90% respectively, and positive likelihood ratios between 8.7 and 14 for different operational trade-offs. Our tool offers an objective, disease-specific, fast (average time $\leq 4$ min) approach to detect fake PTSD, and if adopted, can ensure that healthcare resources and disability concessions reach those genuinely in need, while helping to maintain integrity of clinical data. Moreover, the ability to identify and help patients who might be malingering due to other mental health conditions, poverty or socio-economic compulsions can improve general health outcomes in disadvantaged communities.**

## INTRODUCTION

Diagnosis of mental health disorders typically rely on structured interviews[5] susceptible to intentional fabrication of symptoms[1,2], often referred to as "malingering". In the specific context of clinically diagnosing Post-Traumatic Stress Disorder (PTSD), subjective stressors, a high degree of similarity in presentation, the realtively easy access to information on how to fake PTSD, and financial incentives from a positive diagnosis in legal procedures or disability claims, is known to hinder accurate assessments. Here, we present an approach based on a generative model to flag malingering suspects rapidly, efficiently and accurately; essentially demonstrating a generative artificial intelligence (AI) framework that, paraphrasing Ekman[6], *"can catch a liar"*.

Clinically, PTSD is an anxiety disorder that can develop after experiencing a traumatic event. In the United States, substantial disability compensation may be available for those with mental disorders, which while being an crucial resource for the truely afflicted, incentivizes malingering[3,4]. Thus, faking PTSD to access medical treatment, commit insurance, personal injury and other frauds, or in an attempt to evade criminal liability and penalties[7–10] is unfortunately not rare. While PTSD is a serious mental health condition associated with substance use disorder, mood disorder, anxiety disorder, personality disorder, increased morbidity, and possibly with increased mortality[11,12], false diagnoses can cause substantial financial drain[13,14] to healthcare systems, divert crucial resources from where they are needed[15], and interfere with study outcomes by introducing inaccuracies in clinical data[16]. Accurate disambiguation of true and fake PTSD is therefore of high importance, especially with sources suggesting that over 20% of personal injury cases, as well as 20% of the Veterans seeking combat compensation could be fabricating their condition[3,17–19]. In addition to curbing financial fraud, ability to flag such incongruities can help address the myriad of factors that can drive malingering behavior, including other mental health conditions, a lack of access to healthcare[20,21] and an inability to seek help arising from poverty and other broad-ranging socio-economic conditions.

Fig. 1: **Conceptual framework.** Using a dataset of responses to a validated structured interview for PTSD diagnoses, along with physician-validated clinical diagnoses, we infer generative models for responses for PTSD patients. In our framework for detecting malingering, we flag responses as those which are highly "surprising" (defined as violating infered cross-dependencies between individual response items) or are too simple (lacking complexity in the response patterns typical of non-malingered responses). The precise "malingering condition" shown above is validated from theoretical considerations as well as field experimental data.

Despite the general difficulty in formulating principles to detect malingering[22], multiple standardized tests[23,24] and validity assessment tools[25] have been proposed, with limited success. These tools typically aim to incorporate patterns observed in diagnostic populations that might disambiguate faked symptoms from real ones, or ask similar or related questions multiple times to verify consistency. However, existing approaches do not target specific disorders, almost always require expert interpretation, are subjective, and by design are unlikely to be effective against a malingerer with psychiatric training (See Table 3). Other strategies with physiological monitoring and linguistic analysis[26–29] cannot be easily adopted in structured or semi-structured interviews.

In VeRITAS we leverage the fact that responses in a structured interview have statistical dependencies arising from the nature of the questions themselves, and are modulated by the trait we are aiming to detect $e.g.,$ PTSD pathology. We operationalize this principle without requiring human-understanding of the specific items being presented to the subject; thus making the approach specific to the disease at hand (PTSD), while being potentially generalizable to other disorders if appropriate training data is available.

Our key finding here is that the subtle cross-dependencies between the interview items are challenging to mimic

TABLE 1: Prolific Dataset Average Characteristics and Success Rates as a Percentage of Representation, at 90% sensitivity

| characteristics | mean Completion Time [s] | mean age [years] | no. of participants | success rate (%) |
|---|---|---|---|---|
| Race: Asian | 188.4 | 34.3 | 24 | 8.3 |
| Race: Black | 329.5 | 38.0 | 33 | 18.2 |
| Race: Mixed | 195.8 | 31.4 | 20 | 5.0 |
| Race: Other | 286.2 | 38.5 | 13 | 0.0 |
| Race: White | 186.5 | 42.7 | 220 | 4.1 |
| Sex: Female | 201.5 | 40.5 | 167 | 4.8 |
| Sex: Male | 213.5 | 41.0 | 141 | 7.1 |
| Residence: United Kingdom | 213.5 | 43.1 | 110 | 4.5 |
| Residence: United States | 202.9 | 39.4 | 200 | 6.5 |

TABLE 2: Minimum performance trade-offs (PPV and NPV are calculated assuming prevalence of 30%)

| specificity | sensitivity | ppv | accuracy | npv | LR+ | LR- |
|---|---|---|---|---|---|---|
| 0.90 | $0.856\pm0.021$ | $0.797\pm0.001$ | $0.887\pm0.006$ | $0.933\pm0.009$ | $8.742\pm1.084$ | $0.158\pm0.023$ |
| 0.94 | $0.833\pm0.023$ | $0.863\pm0.002$ | $0.907\pm0.007$ | $0.926\pm0.010$ | $14.11\pm0.337$ | $0.177\pm0.025$ |

on-the-fly, even if the subject is knowledgeable about how such patients tend to respond $i.e.$ with training in the mental health services. Thus, VeRITAS offers a robust approach to verifying response validity, can target specific disorders, can be administered in less that 4 minutes, and require little subjective interpretation.

# RESULTS

## Participants and Data Sources

Our first dataset (referred to as the VA dataset) comprises $n = 304$ participants recruited at a Veterans Health Administration facility for an earlier study[30]. Veterans between the ages of 18 and 89 years were recruited with written informed consent. Once eligibility was determined by the study team, participants completed a PTSD-symptom questionnaire from the CAT-PTSD item bank[30,31], comprising 211 items, including some items from the PTSD Checklist (PCL-5). Participants were also interviewed using the Clinician-Administered Scale for PTSD for Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition (DSM-5), which resulted in 86 PTSD diagnoses, and 218 participants deemed as not having PTSD. We used 60% of the data for inferring our generative models, and the rest was used for validation, including determining the out-of-sample AUC for identifying PTSD vs no-PTSD cases. The possibility of malingering was not recorded in this dataset.

The second dataset generated in the course of the present study is from online surveys conducted by a third party vendor (Prolific) hired by the study team. Referred to here as the Prolific dataset, it comprises $n = 310$ participants (200 in the US, and 110 in the UK), screened for the absence of past or present mental health diagnoses. The participants were asked to fake symptoms of PTSD while taking the VeRITAS interview.

In total, this study considered $n = 624$ participants comprising US Veterans, and drawn from the general population in the US and the UK, with broad representation across sexes and ethnicities. Average completion times and detailed demographic composition of the respondents is shown in Table 1.

## Principle of Characterizing Malingered Responses

Amongst those who have a naive PTSD diagnosis, we aim to flag subjects with a quantifiable high estimated likelihood of malingering. Note that the in VA data set we only have the naive diagnoses information; thus some participants diagnosed with PTSD in the dataset could have been malingering, but are not flagged to be so. To flag malingering sunjects, our first insight is that such response sequences have high average "surprise", $i.e.$ deviate more on average from the context-specific model predictions of item responses. Our other insight is that malingering tend to generate less complex patterns in the response sequences. We understand complexity in the formal sense: more complex objects are less compressible, and random sequences are not significantly compressible. Thus a sequence of all 1's is very compressible and therefore non-random, but a sequence generated from teh sequential tosses of a fair coin is not very compressible, and is highly random. We hypothesize that true response sequences will tend to have maximal randomness, conditional on being constrained by the emergent cross-dependencies. Thus, true responses should have the just right degree of identifiable structure, and no more. More structure will make it less random, and less structure will increase surprise, and thus too much and too little structure are both indicative of malingering. Both of these criteria are based on statistical properties of the response sequence, and do not require a human understanding of the items. Quantitatively,
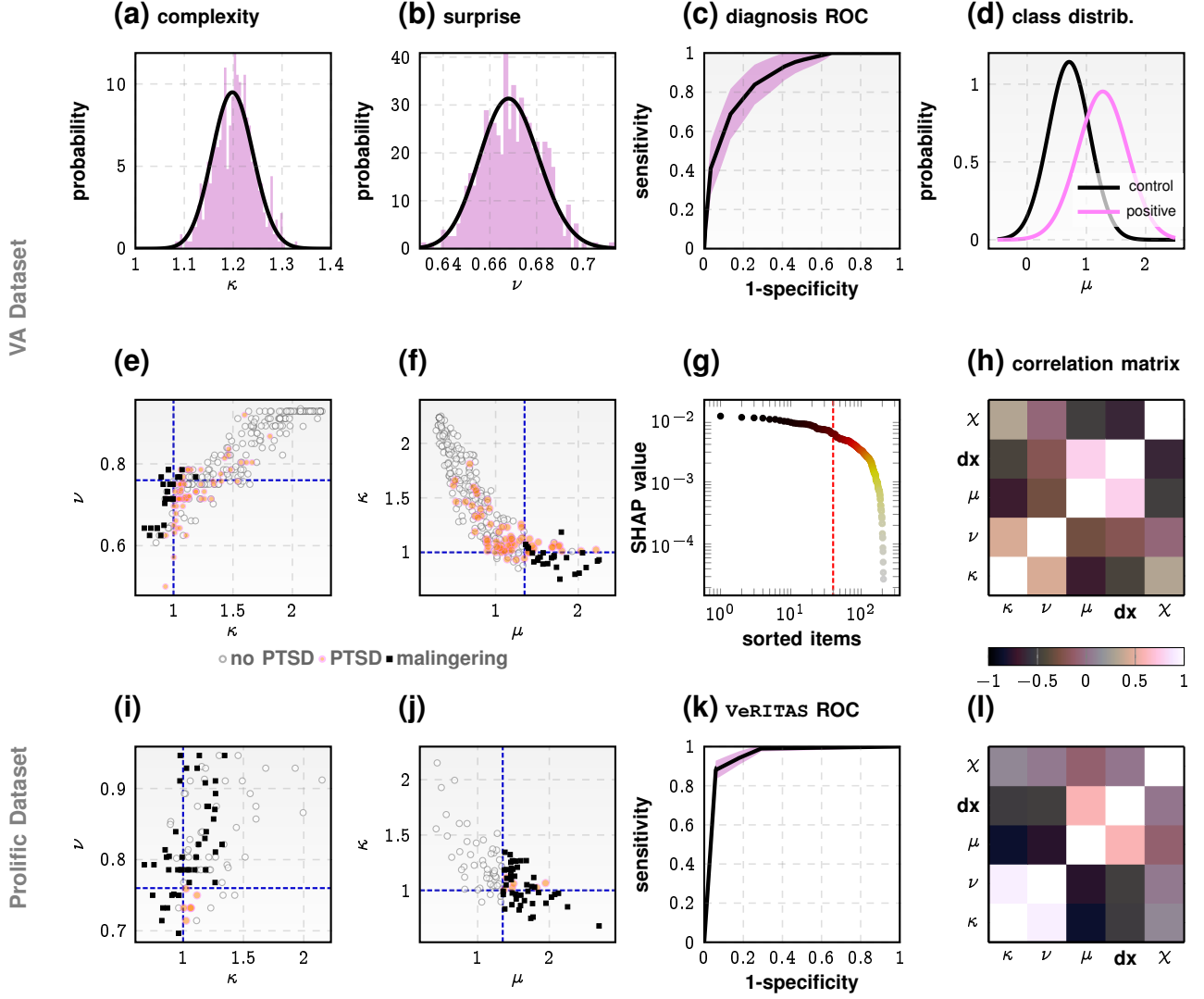
**(a)** complexity    **(b)** surprise    **(c)** diagnosis ROC    **(d)** class distrib.

**(e)**    **(f)**    **(g)**    **(h)** correlation matrix

○ no PTSD ● PTSD ■ malingering

**(i)**    **(j)**    **(k)** VeRITAS ROC    **(l)**

Fig. 2:

for each response sequence x, we compute two quantities: $\kappa(x)$ and $\nu(x)$, referred to as the complexity and the surprise parameters respectively.

The distributions for $\kappa, \nu, \mu$ may be characterized from the VA data, despite the fact that this dataset does not have designation of "malingering". This is because that these quantities are computable from just the response sequences. Concretely, we propose a response sequence x should be flagged as an instance of malingering if for suitably chosen thresholds $\mu_0, \nu_0, \kappa_0$:

$$\chi(x) \triangleq \left(\mu(x) \geqq \mu_0\right) \bigwedge \left(\left(\kappa(x) \leqq \kappa_0\right) \vee \left(\nu(x) \geqq \nu_0\right)\right) \tag{1}$$

This may be paraphrased as "a response sequence has high likelihood of malingering if it 1) produces a naive diagnosis of PTSD with high probability, and 2) is either too surprising or too simple." The decision thresholds are obtained from theoretical considerations and the VA data (See Methods), which allows us to choose thresholds as reflecting specifcity-sensitivity trade-offs (Fig. 2 a,b).

**Integrated Diagnostic Capability**

Given a sequence of responses to a diagnostic interview, our first task is to determine if a particular subject should be diagnosed with PTSD, if the possibility of malingering is ignored. We call this a "naive diagnosis". This diagnostic information might be available to VeRITAS externally (*e.g.* from a physician's assessment of the patients). However, VeRITAS also has an integrated capability for naive diagnosis: we identify separate generative models for 1) the diagnosed set of patients in the training set ($M^+$, some of them might be malingering) and 2) for the patients identified as not having PTSD ($M^0$). Then given the sequence of responses from a new subject, we estimate if that sequence is more likely to have been generated by the model for $M^+$ vs that for $M^0$. We validate this diagnostic capability using the training data described above, and we achieve good disambiguation between $M^+$ and $M^0$, with out-of-sample AUC 0.867 ± 0.008 (See Fig. 2 c), which is at par or

4

better compared to reported tools, $e.g.$, CAT-PTSD[32]. In the `VeRITAS` algorithm, for a given response sequence $\mathbf{x}$, the naive diagnosis risk score is denoted as $\mu(x)$, and referred to as simply the "score" (See Methods for definition and computational approach). Fig. 2,d shows the class specific distributions estimated for $\mu$ for the VA data.

## `VeRITAS` Validation Strategy

Since the VA data does not indicate presence or absence of malingering, we adopted a non-standard approach for validating `VeRITAS`. We assume that the Prolific participants do not have PTSD (based on their screening of not having past or present mental health diagnosis, little or no axiety severity, and being informed to not take the test if experiencing PTSD symtoms), and that all of them were attempting to malinger as directed. This allows us to measure false negative rates, as function of the `VeRITAS` parameters. The succes rates of these participants in getting a "diagnosis" at a sensitivity of 94.2% ($\kappa_0 = 1, \nu_0 = 0.76, \mu_0 = 1.35$) is shown in Table 1. Then we check how many of the VA participants are flagged as malingering amongst the ones with a naive PTSD diagnosis. To determine a lower bound on performance we can assume all of these subjects are false positives (which is unlikely, but nevertheless is an upper bound on false positives as function of `VeRITAS` parameters). This allows us to construct a lower envelop of the ROC curve, and hence estimate a lower bound of the AUC for determining malingering, establishing a minimum AUC of $0.95 \pm 0.02$ at 95% confidence. Two high performance operating point choices, reflecting specificity/sensitivity trade-offs is shown in Table 2, which illustrates that we can achieve $> 83.3\%$ sensitivity with $> 94\%$ specificity, along with $> 86.3\%$ PPV and $> 92.6\%$ NPV, and positive likelihood ratio $> 14.11$. We not eagain that these numbers represent lower bounds on `VeRITAS` performance. The variation of the complexity and surprise parameters for the VA data, along with a set of decision thresholds, is shown in Fig. 2 e,f and for the Prolific data in panels i,j. The correlation matrix between $\kappa, \nu, \mu, dx, \chi$ (where dx is the variable for naive diagnosis) for the VA dataset (panel h) shows that the complexity and teh surprise parameters are well-correlated, and the dx and the $\mu$ parameters are well-correlated, and the malingering decision $\chi$ is not very well correlated to either of these. This pattern is closely replicated in the Prolific dataset (panel l), where dx refers to predicted naive diagnosis. The estimated lower envelop for the ROC curve for the `VeRITAS` algorithm is shown in panel k.

## Item Presentation and Response Sequence in Administration

In `VeRITAS` we present a small random subset of the full item bank to the subjects. Although the models for cross-dependency are inferred using the complete response sequence to all 211 items available, presenting a smaller item set makes the test feasible, while also producing a vast number of possible variations. In the version we validated, 20 items are presented. These are randomly generated from a distribution reflecting the impact of the items in the inferred models on their ability to disambiguate between responses from diagnosed PTSD patients and the control group. Impact is measured using standard SHAP analysis of the inferred models (See Methods). SHAP-values for items are normalized to a probability mass function (pmf), which is then sampled to produce the item set. We can select items adaptively similar to CAT-PTSD[30]; such complications have been presently avoided. Panel g in Fig. 2 shows the sorted estimated SHAP values for the items, and shows the threshold of selection for presentation.

## Q-net Inference and Cross-talk Modeling

The generative models for response sequences inferred separately from the PTSD-positive and PTSD-negative patients are referred to as "Q-nets". Structurally, an Q-net comprises an interdependent collection of local predictors, each aiming to predict the response to a specific item using as features the responses to ther items from the item bank. (Fig. 1). Thus, an Q-net comprises almost as many such position-specific predictors as the length of the response sequence. These individual predictors are implemented as conditional inference trees[33], in which nodal splits have a minimum pre-specified significance in differentiating the child nodes. Thus, each predictor yields an estimated conditional response distribution for each item. The set of items acting as features in each predictor are automatically identified, $e.g.$, in the fragment of PTSD-positive Q-net (Fig **??**), the predictor for item ptsd119 includes the response to item ptsd19, that for ptsd19 includes item ptsd43, that for ptsd43 includes ptsd208, which in turn includes ptsd43, revelaing complex possibly cyclic dependencies. Owing to this recursive expansion, a complete Q-net captures the complexity of the rules guiding the subtle cross-talk emerging between different items in a survey, and can be trained wihout linguistic understanding of the items themselves. The Q-nets inferred for $M^+$ and $M^0$ have 5,852 and 10,301 independent parameters, encoding the probability distributions on the tree leaves, which were identified automatically in the course of inferring the sets of conditional inference trees.

## Time-to-completion and Malingering

In the Prolific data, the average completion time was 3.4 minutes ([2.9, 3.9], 95% CI). We find that malingering flags have no significant correlation with the time to completion. The success rate of defating `VeRITAS` is

observed to have a correlation of $0.48([-.27, .87], 95\%CI)$. Thus, fast responses that might reflect poor quality survey in general[34–36] might not be effective in detecting malingering, espeially if participants are putting more thought into how to fake their symptoms.

### Number of Distinct Variations of `VeRITAS` Implementation

With $r = 20$ items presented, randoml;y chosen from the top $N = 40$ shortlisted from the master item bank (sorted according to SHAP values, as described before), and each item having $L = 5$ possible responses, we can have $\binom{N}{r} = 137,846,528,820$ variations of the test, with $L^r\binom{N}{r} = 13,146,069,414,138,793,945,312,500$ or approximately $10^{25}$ possible responses, which is approximately equal to the number of stars in the observable universe. Thus, it is non-trivial for human subjects to "learn" or "train" to defeat the algorithm.

### Comparison Against State of Art

Structured Interview of Reported Symptoms, 2nd Edition (SIRS-2) has a reported performance of sensitivity of .80, a specificity of .975, and positive and negative predictive powers of more than .90 (based on a base rate of 31.8%), tales 30-40 minutes to complete, and needs expert interpretation and is not disease-specific. In contrast, `VeRITAS` may eb complete in under 4 minutes, can have sensitivity and specificity bith above 90%, has PPV over 86% and NPV over 90% in selected operating points, can be tuned to specific disorders, and may be administered automatically. The crucial difference in `VeRITAS` is the near-impossibility of training to defat it, and its effectiveness in the scenario that the subject has psychiatric training. The principles on which existing tool ssuch as SIRS-2 are based makes it highky unlikely to be effective if the subject is familiar with symptomologies of basic mental disorders, and the approaches generally employed to flag malingering.

# DISCUSSION

In this study, we introduce a novel algorithmic technique, `VeRITAS`, aimed at identifying deception in structured interviews for the diagnosis of Post-Traumatic Stress Disorder (PTSD), a context where malingering presents a significant challenge. The prevalence of malingering, as reported in various studies, underscores the importance of developing reliable, rigorous, and principled methods for detecting deception, with estimates of its prevalence in psychiatric and criminal justice settings ranging widely from 8% to 64%[37–39]. This variability highlights the complexity of the issue and the need for advanced detection techniques[40].

Our approach represents a significant departure from traditional methods of malingering detection, which often rely on domain-specific knowledge[41] and standardized tests developed through extensive research on both genuine patients and malingerers (Table 3). While these traditional methods, including the Structured Interview of Reported Symptoms (SIRS)[23], the Structured Inventory of Malingered Symptomology (SIMS)[24], and validity scales associated with the Minnesota Multiphasic Personality Inventory-2[25], have shown accuracy rates in the range of 85% to 95% depending on the specific problem context[42–45], they come with limitations. These include the requirement for substantial expertise to develop and administer, vulnerabilities to clinician bias, and the potential for false positives and negatives. Furthermore, they may not easily extend to other contexts or disorders, and their effectiveness can be compromised by coaching or prior knowledge of psychiatric symptomology.

Contrastingly, `VeRITAS` leverages a model-based, data-driven approach that identifies structural differences in an individual's response patterns compared to a set of baseline responses, aiming to capture the complex, often non-obvious dependencies between interview questions. This approach utilizes a novel architecture, the Q-net, a nonparametric generative model that maps the inter-question relations and dependencies, enabling the detection of inconsistencies indicative of malingering. This method's efficacy is not solely predicated on the content of responses but on their statistical and structural properties, offering a robust alternative to existing techniques.

The motivation behind fabricating PTSD symptoms is multifaceted, often driven by financial incentives such as disability compensation, insurance claims, or legal benefits. The potentially high prevalence of malingering in some populations not only strains healthcare and legal systems but also undermines the integrity of clinical diagnoses and research, leading to misallocation of resources and potentially hindering the care of individuals genuinely afflicted by PTSD.

Our findings indicate that `VeRITAS` can achieve high sensitivity and specificity in detecting malingering, with performance metrics potentially surpassing or at least comparable to those of existing state-of-the-art techniques, but with several advantages. `VeRITAS` requires less time for administration, does not necessitate domain-specific expertise for interpretation, and minimizes the risk of bias. Moreover, its design makes it challenging for individuals, even those with psychiatric training, to defeat the system through coaching or preparation.

In conclusion, while malingering presents a persistent challenge in the accurate diagnosis and treatment of PTSD, our study offers a promising new direction for detection methods. By employing a sophisticated, data-driven approach that transcends the limitations of traditional methods, `VeRITAS` provides a powerful tool for

clinicians and researchers. Its adoption could significantly enhance the integrity of PTSD diagnoses, ensuring that resources are allocated to those genuinely in need and supporting the broader goals of psychiatric care and research. However, further research and validation across diverse populations and settings are essential to fully realize its potential and applicability in clinical practice.

# METHODS

## 1. DEFINITIONS & NOTATION

**Definition 1** (Survey). *A survey for the purpose of this work is a structured interview, consisting of a finite number of questions (items) posed to a set of participants, with these items drawn from a finite item bank, and whose responses must be one froma pre-specified set fo choices, $e.g.$, the Likert scale, with missing values for the responses allowed.*

**Definition 2** (Response vector). *A response vector is the set of responses to a survey from a single participant, typically assuming that not all items are posed, and allows for the possibility that some responses are missing.*

A Q-net, as described here, is a model of the response dependency structure for questions (items) posed to participants in a survey. The Q-net explicitly estimates individual conditional distributions of each item response, which collectively serve as a model of the full joint distribution of the responses.

**Definition 3** (Q-net). *Let $X \sim P$ be an $n$-dimensional discrete random vector supported on a finite set $\Sigma$ and following distribution $P$, i.e.*

$$X = (X_1, \ldots, X_n) \sim P, \qquad \mathrm{supp}(X) = \Sigma = \prod_{i=1}^{n} \Sigma_i \ \ \text{with } |\Sigma| < \infty.$$

*For $i = 1, \ldots, n$, let $P_i := P(X_i \,|\, X_j = x_j \text{ for } j \neq i)$ denote the conditional distribution of $X_i$ given the values of the other components of $X$. Finally, for each $i = 1, \ldots, n$, let $\Phi_i^P$ denote an estimate of the distribution $P_i$. Then the set $\Phi^P := \{\Phi_i^P\}_{i=1}^{n}$ is called a* Quasinet (Q-net). *Identifying the true distribution $P$ as the one describing the joint statistics of the responses from a survey with $n$ items, we also refer to $\Phi^P$ as the Q-net for the survey $P$.*

When $P$ is clear from context, we may omit the superscript and simply write $\Phi = \{\Phi_i\}$ to denote the Q-net. The motivation for Definition 3 is that the collection of all estimators $\Phi = \{\Phi_i\}$ contained in a Q-net represents the set of all inferred dependencies from the observed ecosystem. While the definition allows for arbitrary method of algorithm to construct the estimators $\Phi_i$, the utility of a Q-net clearly depends primarily on the properties of the $\Phi_i$. In this study, we aim to minimize the set of a priori assumptions on the overall model structure to allow the complex dependencies present in $P$ to emerge. To that end, throughout this work all Q-nets are computed using conditional inference trees[46] (a variant of classification and regression trees) to compute each $\Phi_i$. In general each Q-net component $\Phi_i$ is computed independently from the other $\Phi_j$, which allows a network structure to emerge amongst these estimators.

An important quantity for an inferred Q-net is the persistence function $\omega_{\mathbf{x}}$.

**Definition 4** (Persistence Function). *Given a survey $P$ inducing the Q-net $\Phi^P$ and a response vector $\mathbf{x} = (x_1, \ldots, x_n)$, the persistence $\omega_{\mathbf{x}}$ of $\mathbf{x}$ in the population modeled by the Q-net:*

$$\omega_{\mathbf{x}}^P := \Pr(\mathbf{x} \in P) = \prod_{i=1}^{n} \Phi_i^P(X_i = x_i \,|\, X_j = x_j, j \neq i) \tag{2}$$

The persistence function $\omega_{\mathbf{x}}^P$, as the name suggests, is the probability that $\mathbf{x}$ persists, $i.e.$, $Pr(\mathbf{x} \to \mathbf{x})$ for the population modeled by the Q-net $P$, with $1 - \omega_{\mathbf{x}}^P$ being the probability that $\mathbf{x}$ is altered by a random perturbation.

We will show that if for two inferred Q-net models $P, Q$, we have $\omega_{\mathbf{x}}^P \geqq \omega_{\mathbf{x}}^Q$, then it is more likely that model $P$ generated $\mathbf{x}$. This is an important result that justifies the definition of the score parameter in Defn. 6.

The Q-net allows us to rigorously compute bounds on the probability of a spontaneous change from one response vector to another, induced by spontaneous chance variations. Not all perturbations in a vector are either likely or contextually meaningful. With an exponentially exploding number of possibilities in which a vector over a large set of items can vary, it is computationally intractable to directly model all possible dependecies; nevertheless, we can constrain the possibilities using the patterns we uncover via the Q-net construction. A key piece of this approach is to design an intrinsic distance between response vectors, which is reflective of this underlying dependency structure.

**Definition 5** (q-distance). *Let $\Phi^P = \{\Phi_i^P\}_{i=1}^{n}$ and $\Phi^Q = \{\Phi_i^Q\}_{i=1}^{n}$ denote Q-nets on populations $P$ and $Q$, and suppose $\mathbf{x} = (x_1, \ldots, x_n)$ and $\mathbf{y} = (y_1, \ldots, y_n)$ are samples of $X \sim P$ and $Y \sim Q$ respectively. Then the*

*q-distance $\theta_{P,Q}(\mathbf{x}, \mathbf{y})$ between $\mathbf{x}$ and $\mathbf{y}$ is*

$$\theta_{P,Q}(\mathbf{x}, \mathbf{y}) := \frac{1}{n} \sum_{i=1}^{n} \left[ \mathbb{J}^{\frac{1}{2}} \left( \Phi_i^P(X_i | X_j = x_j, j \neq i) \| \Phi_i^Q(Y_i | Y_j = y_j, j \neq i) \right) \right]$$

*where $\mathbb{J}$ denotes the Jensen-Shannon divergence[47].*

For brevity, we may write simply $\theta$ (dropping the suffixes) if the populations are clear from context. Since the Jensen-Shannon distance $\mathbb{J}$ is a legitimate metric[48] on the set of probability distributions (unlike KL-divergence), $\theta$ inherits nonnegativity, symmetry, and respects the triangle inequality; it follows that q-distance is a (pseudo)-metric on $\Sigma$. Note that, being a pseudo-metric implies that we may have $\theta(\mathbf{x}, \mathbf{y}) = 0$ for $\mathbf{x} \neq \mathbf{y}$, i.e. distinct vectors can induce the same distributions over each index, and thus have zero distance. This is in fact desirable, since we do not want our distance to be sensitive to changes that are not meaningful. The intuition is that not all variations are equally important or likely. Moreover, we show in Theorem 1 that the log-likelihood of a vector $\mathbf{x}$ transitioning to $\mathbf{y}$ scales with $\theta(\mathbf{x}, \mathbf{y})$, allowing us to directly estimate the probability of spontaneous (or sequential) jumps between abundance profiles.

**Theorem 1** (Probability Bound). *Given a vector $\mathbf{x}$ of length $n$ from $P$ that transitions to $\mathbf{y}$ from $Q$, we have the following bounds at significance level $\alpha$.*

$$\omega_{\mathbf{y}} e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta(\mathbf{x},\mathbf{y})} \geqq Pr(\mathbf{x} \to \mathbf{y}) \geqq \omega_{\mathbf{y}} e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta(\mathbf{x},\mathbf{y})} \tag{3}$$

*where $\omega_{\mathbf{y}}$ is the persistence of $\mathbf{y}$ (Def. 4), and $\theta(\mathbf{x}, \mathbf{y})$ is the q-distance between $\mathbf{x}, \mathbf{y}$ (Def. 5).*

*Proof.* See later in Section 3. $\square$

Theorem 1 gives theoretical backing to the claim that samples generated by the Q-net indeed reflect likely perturbation possibilities from the current state. Thus we can use the Q-net to draw contextually realistic samples that respect the cross dependencies and reduce surprise (that is, the Q-net-inferred conditional distributions can be used to generate approximate samples from the population $P$).

**Remark 1** (Neighborhood Structure). *It follows from Th. 1 that we have for some constant $C$,*

$$\ln \left| \frac{Pr(x \to y)}{Pr(y \to y)} \right| \leqq C\theta(x, y) \tag{4}$$

*implying for all response vectors $y$ within a small neighborhood of $x$ (small in metric $\theta$), we have:*

$$Pr(y \to x) \approx Pr(x \to x) \tag{5}$$

*which reveals an important special structure on local neighborhoods.*

## 2. VERITAS ANALYSIS

**Definition 6** (Algorithm `VeRITAS` Parameters). *We introduce three parameters referred tro as teh complexity, surprise and score parameters ($\kappa, \nu, \mu$ respectively) for a given response vector $x$:*

$$\text{complexity: } \kappa \triangleq -\frac{1}{|x|} \ln Pr(x \to x | M^+) = -\frac{\ln \omega_x^{M^+}}{|x|} \tag{6}$$

$$\text{surprise: } \nu \triangleq \mathbf{E}_i \left( 1 - \Phi_i^{M^+}(x_{-i}) |_{x_i} \right) \tag{7}$$

$$\text{score: } \mu \triangleq \frac{\ln Pr(x \to x | M^+)}{\ln Pr(x \to x | M^0)} = \frac{\ln \omega_x^{M^+}}{\ln \omega_x^{M^0}} \tag{8}$$

*where $M^+$ indicates the sub-population exhibiting a particular trait of interest e.g. a mental health disorder such as PTSD, and $M^0$ is the control sub-population where this trait is absent.*

**Definition 7** (Malingering property). *A response vector $x$ is defined to have the malingering property if:*

$$\chi(x) \triangleq \left( \mu(x) \geqq \mu_0 \right) \bigwedge \left( \left( \kappa(x) \leqq \kappa_0 \right) \vee \left( \nu(x) \geqq \nu_0 \right) \right) \tag{9}$$

The decision thresholds $\kappa_0, \nu_0, \mu_0$ are inferred from survey data.

**Lemma 1** (Complexity). *For a survey with $n$ items, and assuming $L$ to the number of possible responses to each item, the unconditional probability of a response vector $x$ occurring among all feasible responses is bounded above by $(e^\kappa/L)^n$, where $\kappa(x)$ is the complexity parameter for response $\mathbf{x}$.*

*Proof.* Let $\kappa(x) \leqq \kappa'$. From Def. 6, we have for a response vector $x$,

$$-\frac{1}{n} \ln \omega_x \leqq \kappa' \Rightarrow \omega_x \geqq e^{-n\kappa'} \tag{10}$$

Summing on both sides over all responses $x$ with $\kappa(x) \leqq \kappa'$ (assume there are $N_x$ such sequences), we have:

$$1 \geqq \sum_x \omega_x \geqq \sum_x e^{-n\kappa'} \tag{11}$$

where the first inequality follows from observing that responses very close to $x$ in the q-distance metric have a specific structure, namely $\omega_x \approx Pr(y \to x)$ (See Remark 1) and responses further away have smaller jump probabilities, which then implies:

$$N_x \sum_x e^{-n\kappa'} \leqq 1 \Rightarrow N_x \leqq e^{n\kappa'} \tag{12}$$

The result then follows from noting that the complete set of possible responses has the size $L^n$. □

Lemma 1 justifies why a low value of $\kappa$ implies the possibility of an un-natural response, because the odds of generating such a response is remarkably small.

**Corollary 1** (Algorithmic Complexity). *, The algorithmic complexity of a response $x$ conditional on the number of survey items $n$ is at most $\kappa(x) + O(1)$.*

*Proof.* This follows from noting that a set of cardinality $L^m$ has a algorithmic complexity of $m + O(1)$, since words of length $m$ are sufficient to encode the index of any element of the set, and thus can be uniquely identified. Sinec we can calculate $\kappa' = \kappa(x)$ for any $x$, and since the set of all $x$ for a given value of $\kappa'$ belongs to a set of size at most $e^{-n\kappa'}$, the result follows. □

**Lemma 2** (Surprise). *For any response vector $x$, we have:*

$$\nu(x) \leqq 1 - e^{-\kappa(x)} \tag{13}$$

*Proof.* Denoting $\Phi_i(x_{-i})|_{x_i}$ as $a_i$, we note that $\omega_x^{1/n}$ is the geometric mean of the vector of $a_i$s, while $\mathbf{E}_i\left(\Phi_i(x_{-i})|_{x_i}\right)$ is the arithmetic mean of the same vector, which then completes the proof by noting:

$$-\mathbf{E}_i\left(\Phi_i(x_{-i})|_{x_i}\right) \leqq -\omega_x^{1/n} \Rightarrow \nu(x) \leqq 1 - \omega_x^{1/n} \tag{14}$$

□

## Interpretation on Why the Defined Property Identifies Malingering

Lemma 2 indicates that the requirement of an upper bound on the surprise and a lower bound on the complexity are both aiming to flag responses which are unlikely to appear when the data (responses) are being generated by the underlying process corresponding to the phenotype of interest (PTSD). When such unlikely responses do appear appear nevertheless, it is likely that they are not being generated by the correct underlying process. One can attempt to fake responses that might seem to increase the odds of a positive diagnosis, but the respondant must replicate the cross-dependencies closely enough (build in enough structure) so that the deviation from the expected responses is limitd (limited surprise requirment). But building in too much structure will reduce the complexity too much (too much structure reduces complexity, since there are fewer highly structured sequences), which will then fail the complexity lower bound.

Note that the remaining condition $\mu(x) \geqq \mu_0$ is a diagnosis criterion for the trait of interest ($M^+$), and may be replaced with a different condition if available for identifying participants with the $M^+$ trait. This particular form follows from a straightforward Bayesian argument on estimating the posterior.

## 3. PROOF OF THEOREM 1

**Theorem 2** (Probability bound). *Given a sequence $x$ of length $N$ that transitions to a strain $y \in Q$, we have the following bounds at significance level $\alpha$.*

$$\omega_y^Q e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta(x,y)} \geqq Pr(x \to y) \geqq \omega_y^Q e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta(x,y)} \tag{15}$$

*where $\omega_y^Q$ is the membership probability of strain $y$ in the target population $Q$ (See Def. **??**), and $\theta(x,y)$ is the q-distance between $x, y$ (See Def. 5).*

*Proof.* Using Sanov's theorem[47] on large deviations, we conclude that the probability of spontaneous jump from strain $x \in P$ to strain $y \in Q$, with the possibility $P \neq Q$, is given by:

$$Pr(x \to y) = \prod_{i=1}^{N}\left(\Phi_i^P(x_{-i})|_{y_i}\right) \tag{16}$$

Writing the factors on the right hand side as:

$$\Phi_i^P(x_{-i})|_{y_i} = \Phi_i^Q(y_{-i})|_{y_i}\left(\frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}}\right) \tag{17}$$

we note that $\Phi_i^P(x_{-i})$, $\Phi_i^Q(y_{-i})$ are distributions on the same index $i$, and hence:

$$|\Phi_i^P(x_{-i})_{y_i} - \Phi_i^Q(y_{-i})_{y_i}| \leqq \sum_{y_i \in \Sigma_i} |\Phi_i^P(x_{-i})_{y_i} - \Phi_i^Q(y_{-i})_{y_i}| \tag{18}$$

Using a standard refinement of Pinsker's inequality[49], and the relationship of Jensen-Shannon divergence with total variation, we get:

$$\theta_i \geqq \frac{1}{8}|\Phi_i^P(x_{-i})_{y_i} - \Phi_i^Q(y_{-i})_{y_i}|^2 \Rightarrow \left|1 - \frac{\Phi_i^Q(y_{-i})_{y_i}}{\Phi_i^P(x_{-i})_{y_i}}\right| \leqq \frac{1}{a_0}\sqrt{8\theta_i} \tag{19}$$

where $a_0$ is the smallest non-zero probability value of generating the entry at any index. We will see that this parameter is related to statistical significance of our bounds. First, we can formulate a lower bound as follows:

$$\log\left(\prod_{i=1}^{N} \frac{\Phi_i^P(x_{-i})_{|y_i}}{\Phi_i^Q(y_{-i})_{|y_i}}\right) = \sum_i \log\left(\frac{\Phi_i^P(x_{-i})_{|y_i}}{\Phi_i^Q(y_{-i})_{|y_i}}\right) \geqq \sum_i \left(1 - \frac{\Phi_i^Q(y_{-i})_{y_i}}{\Phi_i^P(x_{-i})_{y_i}}\right) \geqq \frac{\sqrt{8}}{a_0}\sum_i \theta_i^{1/2} = -\frac{\sqrt{8}N}{a_0}\theta \tag{20}$$

Similarly, the upper bound may be derived as:

$$\log\left(\prod_{i=1}^{N} \frac{\Phi_i^P(x_{-i})_{|y_i}}{\Phi_i^Q(y_{-i})_{|y_i}}\right) = \sum_i \log\left(\frac{\Phi_i^P(x_{-i})_{|y_i}}{\Phi_i^Q(y_{-i})_{|y_i}}\right) \leqq \sum_i \left(\frac{\Phi_i^Q(y_{-i})_{y_i}}{\Phi_i^P(x_{-i})_{y_i}} - 1\right) \leqq \frac{\sqrt{8}N}{a_0}\theta \tag{21}$$

Combining Eqs. 20 and 21, we conclude:

$$\omega_y^Q e^{\frac{\sqrt{8}N}{a_0}\theta} \geqq Pr(x \to y) \geqq \omega_y^Q e^{-\frac{\sqrt{8}N}{a_0}\theta} \tag{22}$$

Now, interpreting $a_0$ as the probability of generating an unlikely event below our desired threshold (*i.e.* a "failure"), we note that the probability of generating at least one such event is given by $1 - (1 - a_0)^N$. Hence if $\alpha$ is the pre-specified significance level, we have for $N >> 1$:

$$a_0 \approx (1 - \alpha)/N \tag{23}$$

Hence, we conclude, that at significance level $\geqq \alpha$, we have the bounds:

$$\omega_y^Q e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta} \geqq Pr(x \to y) \geqq \omega_y^Q e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta} \tag{24}$$

$\square$

TABLE 3: Summary of malingering/deception detection methods.

| Approach/Tool | Description | Noted Accuracy/Performance |
|---|---|---|
| Structures Interview of Reported Symptoms (SIRS)[23] | An interview-based measure with multiple detection strategies. | No specific accuracy rate mentioned, but noted as a robust instrument. |
| Structured Inventory of Malingered Symptomology (SIMS)[24] | A paper-and-pencil screening device for detecting malingering with a sensitivity for detecting malingering of 95.6% in a study with college students. | 95.6% sensitivity in a specific study context. |
| Minnesota Multiphasic Personality Inventory-2 (MMPI-2)[25] | A self-report measure assessing personality and psychopathology. Certain validity scales were developed to uncover malingering. | Not specified, but noted flaws and potential for false positives. |
| Millon Clinical Multiaxial Inventory MCMI-III | A self-report scale focusing on personality disorders. | No specific accuracy rate mentioned. |
| Miller Forensic Assessment of Symptoms (M-FAST) | A brief screening measure for malingered mental illness in forensic settings. | Noted issues with low internal consistency on some scales. |
| Human Lie Detectors | Not a specific tool, but the general method of humans attempting to discern lies from truth. | People are generally poor lie detectors. |
| Arousal-Based Approaches (like the Polygraph) | Techniques that rely on physiological responses. | Criticized for poor validity and high rate of false positives. |

| | | |
|---|---|---|
| Cognitive Load-Inducing Approaches | Techniques that view deception as a cognitive act that generally imposes greater cognitive load on respondents than honesty does. | No specific accuracy rate mentioned. |
| Autobiographical Implicit Association Test (aIAT)[43] | Designed to determine whether respondents possess actual autobiographical memories. | 91% accuracy rate in identifying genuine autobiographical memories. |
| Timed Antagonistic Response Alethiometer (TARA)[44] | A computer-administered, response time-based method of lie detection. | 85% accuracy rate. |
| Detecting Faked Identities With Unexpected Questions and Mouse Movements | Technique using computer mouse movements in conjunction with unexpected questions to uncover faked identities[45] | 95% accuracy rate. |
| Time-Restricted Integrity Confirmation (TRI-Con) | A cognitive load-inducing technique with potential to uncover different kinds of deception including malingering. | Up to 89% accuracy rate. |
| Activation-Decision-Construction-Action Theory (ADCAT) | A theory of high-stakes deception. | No specific accuracy rate mentioned, but this is more of a theoretical foundation rather than a specific tool or method. |

# REFERENCES

[1] Rogers, R. An overview of malingering and its assessment. *Psychiatric Clinics of North America* **20**, 15–27 (1997).

[2] Rogers, R. E. *Clinical assessment of malingering and deception* (Guilford Press, 2008).

[3] Frueh, B., Grubaugh, A., Elhai, J. & Buckley, T. Us department of veterans affairs disability policies for posttraumatic stress disorder: administrative trends and implications for treatment, rehabilitation, and research. *Am J Public Health* **97**, 2143–5 (2007).

[4] Taylor, S., Frueh, B. C. & Asmundson, G. J. G. Detection and management of malingering in people presenting for treatment of posttraumatic stress disorder: Methods, obstacles, and recommendations. *Journal of Anxiety Disorders* **21**, 22–41 (2007).

[5] Ali, S., Jabeen, S. & Alam, F. Multimodal approach to identifying malingered posttraumatic stress disorder: A review. *Innovations in Clinical Neuroscience* **12**, 12 (2015).

[6] Ekman, P. & O'Sullivan, M. Who can catch a liar? *American Psychologist* **46**, 913–920 (1991).

[7] Guriel, J. & Fremouw, W. Assessing malingered posttraumatic stress disorder: A critical review. *Clinical Psychology Review* **23**, 881–904 (2003).

[8] Salloway, S., Southwick, S. & Sadowsky, M. Opiate withdrawal presenting as posttraumatic stress disorder. *Hospital and Community Psychiatry* **41**, 666–667 (1990).

[9] Resnick, P. J., West, S. & Payne, J. W. Malingering of posttraumatic disorders. In Rogers, R. (ed.) *Clinical assessment of malingering and deception*, 109–127 (Guilford Press, 2008), 3 edn.

[10] Burkett, B. G. & Whitley, G. *Stolen valor: How the Vietnam generation was robbed of its heroes and history* (Verity Press, 1998).

[11] Goldstein, R. B. *et al.* The epidemiology of dsm-5 posttraumatic stress disorder in the united states: results from the national epidemiologic survey on alcohol and related conditions-iii. *Social psychiatry and psychiatric epidemiology* **51**, 1137–1148 (2016).

[12] Schnurr, P. P., Lunney, C. A., Bovin, M. J. & Marx, B. P. Posttraumatic stress disorder and quality of life: Extension of findings to veterans of the wars in iraq and afghanistan. *Clinical psychology review* **29**, 727–735 (2009).

[13] LoPiccolo, C., Goodkin, K. & Baldewicz, T. Current issues in the diagnosis and management of malingering. *Ann Med* **31**, 166–174 (1999).

[14] Oboler, S. Disability evaluations under the department of veterans affairs. In Rondinelli, R. & Katz, R. (eds.) *Impairment Rating and Disability Evaluation*, 187–217 (W. B. Saunders, Philadelphia, PA, 2000).

[15] Taylor, S. *Clinician's Guide to Treating PTSD: A Cognitive-Behavioral Approach* (Guilford Press, New York, 2006).

[16] Rosen, G. Dsm's cautionary guideline to rule out malingering can protect the ptsd data base. *J Anxiety Disorders* **20**, 530–535 (2006).

[17] Marx, B. & Holowka, D. Ptsd disability assessment. *PTSD Res Q* **22**, 1–6 (2011).

[18] Rogers, R., Sewell, K. & Goldstein, A. Explanatory models of malingering: a prototypical analysis. *Law & Hum Behav* **18**, 543–52 (1994).

[19] Lees-Haley, P. Mmpi-2 base rates for 492 personal injury plaintiffs: implications and challenges for forensic assessment. *J Clin Psychol* **53**, 745–55 (1997).

[20] Park, L., Costello, S., Li, J., Lee, R. & Jacobson, K. C. Race, health, and socioeconomic disparities associated with malingering in psychiatric patients at an urban emergency department. *General Hospital Psychiatry* **71**, 121–127 (2021).

[21] Muntaner, C., Eaton, W. W., Miech, R. & O'campo, P. Socioeconomic position and major mental disorders. *Epidemiologic reviews* **26**, 53–62 (2004).

[22] Drob, S. L., Meehan, K. B. & Waxman, S. E. Clinical and conceptual problems in the attribution of malingering in forensic evaluations. *The journal of the American Academy of Psychiatry and the Law* **37**, 98–106 (2009).

[23] Wong, S. & O'Sullivan, M. The structured interview of reported symptoms (sirs): An overview. *Assessment* **12**, 289–307 (2005).

[24] Smith, G. P. & Burger, G. K. Detection of malingering: validation of the structured inventory of malingered symptomatology (sims). *Journal of the American Academy of Psychiatry and the Law Online* **25**, 183–189 (1997).

[25] Ben-Porath, Y. S. *Interpreting the mmpi-2-rf* (U of Minnesota Press, 2012).

[26] Ekman, P. & O'Sullivan, M. Who can catch a liar? *American Psychologist* **46**, 913 (1991).

[27] Mihalcea, R. & Strapparava, C. The lie detector: Explorations in the automatic recognition of deceptive language. In *ACL-IJCNLP 2009* (2009).

[28] Burgoon, J. K., Blair, J. P. & Strom, R. E. Cognitive biases and nonverbal cue availability in detecting deception. *Human Communication Research* **34**, 572–599 (2008).

[29] Zhou, L., Burgoon, J. K., Nunamaker Jr, J. F. & Twitchell, D. Automating linguistics-based cues for detecting deception in text-based asynchronous computer-mediated communications. *Group Decision and Negotiation* **13**, 81–106 (2004).

[30] Brenner, L. A. *et al.* Development and validation of computerized adaptive assessment tools for the measurement of posttraumatic stress disorder among us military veterans. *JAMA Network Open* **4**, e2115707–e2115707 (2021).

[31] Gibbons, R. D. *et al.* Development of a computerized adaptive test for depression. *Archives of general psychiatry* **69**, 1104–1112 (2012).

[32] Brenner, L. A. *et al.* Development and validation of computerized adaptive assessment tools for the measurement of posttraumatic stress disorder among us military veterans. *JAMA Network Open* **4**, e2115707 (2021). URL https://doi.org/10.1001/jamanetworkopen.2021.15707.

[33] Hothorn, T., Hornik, K. & Zeileis, A. Unbiased recursive partitioning: A conditional inference framework. *JOURNAL OF COMPUTATIONAL AND GRAPHICAL STATISTICS* **15**, 651–674 (2006).

[34] Tourangeau, R., Rips, L. J. & Rasinski, K. The psychology of survey response (2000).

[35] Malhotra, N. Completion time and response order effects in web surveys. *Public opinion quarterly* **72**, 914–934 (2008).

[36] Callegaro, M., Yang, Y., Bhola, D. S., Dillman, D. A. & Chin, T.-Y. Response latency as an indicator of optimizing in online questionnaires. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique* **103**, 5–25 (2009).

[37] McDermott, B., Dualan, I. & Scott, C. Malingering in the correctional system: does incentive affect prevalence? *Int'l J L & Psychiatry* **36**, 287–92 (2013).

[38] Schmidt, T., Krüger, M. & Ullmann, U. Base rate of probable malingering and its indicators in the assessment of mental disorders-retrospective analysis of a sample of forensic psychological evaluations. *Die Rehabilitation* **59**, 231–236 (2020).

[39] Matto, M., McNiel, D. E. & Binder, R. L. A systematic approach to the detection of false ptsd. *The journal of the American Academy of Psychiatry and the Law* **47**, 325–334 (2019).

[40] DePaulo, B. M., Lindsay, J. J., Malone, B. E., Muhlenbruck, L. & Charlton, K. Cues to deception. *Psychological Bulletin* **129**, 74–118 (2003).

[41] Walczyk, J. J., Sewell, N. & DiBenedetto, M. B. A review of approaches to detecting malingering in forensic contexts and promising cognitive load-inducing lie detection techniques. *Frontiers in psychiatry* **9**, 700 (2018).

[42] Rogers, R. & Correa, A. A. Determinations of malingering: Evolution from case-based methods to detection strategies. *Psychiatry, Psychology and Law* **15**, 213–223 (2008).

[43] Sartori, G., Agosta, S., Zogmaister, C., Ferrara, S. D. & Castiello, U. How to accurately detect autobiographical events. *Psychological science* **19**, 772–780 (2008).

[44] Gregg, A. P. When vying reveals lying: The timed antagonistic response alethiometer. *Applied Cognitive Psychology: The Official Journal of the Society for Applied Research in Memory and Cognition* **21**, 621–647 (2007).

[45] Monaro, M., Gamberini, L. & Sartori, G. The detection of faked identity using unexpected questions and mouse dynamics. *PloS one* **12**, e0177851 (2017).

[46] Sarda-Espinosa, A., Subbiah, S. & Bartz-Beielstein, T. Conditional inference trees for knowledge extraction from motor health condition data. *Engineering Applications of Artificial Intelligence* **62**, 26–37 (2017).

[47] Cover, T. M. & Thomas, J. A. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)* (Wiley-Interscience, New York, NY, USA, 2006).

[48] Fuglede, B. & Topsoe, F. Jensen-shannon divergence and hilbert space embedding. In *International Symposium onInformation Theory, 2004. ISIT 2004. Proceedings.*, 31 (IEEE, 2004).

[49] Fedotov, A. A., Harremoës, P. & Topsoe, F. Refinements of pinsker's inequality. *IEEE Transactions on Information Theory* **49**, 1491–1498 (2003).

## SOFTWARE AVAILABILITY

Software for inferring Q-nets is available as an open-source python package `quasinet`, and can be installed from the standard Python code registry.

## ACKOWLEDGEMENTS, AUTHOR CONTRIBUTIONS AND REGULATORY APPROVALS